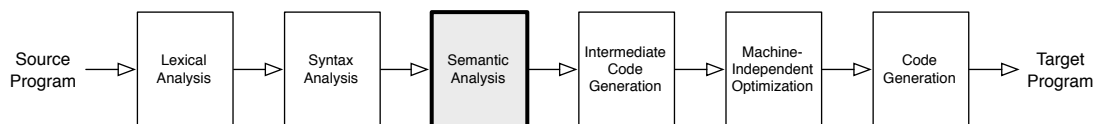


Plan



This week we discuss the design of a type-checker, as well as subtyping and object types. The first set of readings supports these topics.

We then consider two applications of type systems for checking *information flow* policies for security and for precision guarantees in approximate programming. These topics are supported by readings of research and survey papers. When reading research papers, aim to understand the big ideas rather than the specific details. We will discuss these papers (and the associated questions) together in class next week (instead of in meetings).

Readings

Type checker, object types, subtyping:

- Review earlier readings from Cardelli (1-4) and EC (4.2).
- *Foundations of Object Oriented Languages*, Kim Bruce, 2002. Chapters 2 – 3. Skim **as needed** for review of object-oriented basics and subtyping. (Will email about access to copies.)
- “Type Systems,” Luca Cardelli, Section 6 (Subtyping), pages 28–30. (Variant types capture the essence of ML *datatypes*, and are related to Scala *case classes* or other languages’ *enumerations*.)

Type system applications in information flow, security, and approximate programming:

- “Language-Based Information-Flow Security,” Andrei Sabelfeld and Andrew C. Myers, *IEEE Journal on Selected Areas in Communications*, 2003. **Sections I – III**.
<http://www.cs.cornell.edu/andru/papers/jsac/sm-jsac03.pdf>
- “EnerJ: Approximate Data Types for Safe and General Low-Power Computing,” Adrian Sampson, *et al.*, *Programming Language Design and Implementation*, 2011. **Sections 1 – 3.1**.
<http://adriansampson.net/media/papers/enerj-pldi2011.pdf>
There is also a general-audience article that is a good companion to the technical paper:
<http://spectrum.ieee.org/computing/software/enerj-the-language-of-goodenough-computing>

Exercises

1. Design the `tc` package for your IC compiler. This package will contain the code to perform the semantic checks outlined in the IC specification. The main class in the package will be a `TypeCheck` class whose job is to take an AST and annotate each `Expr` node with a type, where `Expr` node is the abstract class from which all expression nodes are derived — your class may be called something different. You will need to extend that class with a `type` field (and accessors) to store the `Type` determined by the typechecker. (Here, I’m using `Type` to refer to the abstract class from which the AST classes representing types are derived — again, your class name may be different.)

Please come to the tutorial meeting with a design detailed enough to discuss the following items:

- Draw the AST for the expression `x + 3 == 7 || a[1] > -x` using your AST package. Annotate each node in the tree with the type corresponding to that expression, given the typing environment `E = int x, int[] a`.
 - The typing rules require you to determine whether one type is a subtype of another. How will you implement subtyping for your `Type` objects?
 - Sketch the implementations of the type checker's code for your AST node class corresponding to each of the following:
 - a unary expression (`!e` or `-e`).
 - an array access
 - a variable access
 - a field access
 - an assignment statement
 - Other than the changes described above, how will you change the `ast` package to support type checking, if at all?
2. Cardelli (Section 1) discusses nominal and structural type equivalence, as do Cooper and Torczon (EC Chapter 4). What is the difference? Which does Java (and IC) use? What about other statically-typed languages you have used (*e.g.*, ML, Scala, Haskell, C#, ...)? Does the distinction matter in dynamically-typed languages (*e.g.*, Racket, Scheme, Python, Ruby, JavaScript, ...)? Do any languages mix the two?

Both authors describe tradeoffs between nominal and structural typing. Do you agree with them? Which is better? Which issues should you worry about?

3. Java's ternary expression `e1 ? e2 : e3` evaluates to `e2` if `e1` is true, and `e3` if `e1` is false.
- (a) Are the following expressions and statements well-typed, *in the sense that no runtime type error could occur as a result of using them in any program?* (Assume `class B extends A`.) For each well-typed ternary expression, what is the most precise type possible? Assume `b` is a boolean variable.
- `b ? 10 : 20`
 - `b ? 10 : true`
 - `A x = (b ? new B() : new A());`
 - `B x = (b ? new A() : new B());`
 - `b ? null : new A()`
- (b) Extend the IC type system to include this construct. Be sure to assign the most precise type possible to a ternary expression, since this will allow the expression to be used in the most contexts (*e.g.*, `b ? null : null` could be given type `A`, `B`, or `Null`, but the third is the most precise type since `A` and `B` are both supertypes of `Null`.) Show the derivation for (iii) to illustrate how the rule works.
4. Java arrays use the following subtyping rule:

$$\frac{A \leq B}{A[] \leq B[]}$$

Demonstrate why this rule causes the type system to be unsound by writing a short program that would cause a run-time type error when executed. (Try running your program. If your answer is correct, you really can make it crash with a runtime type error.)

5. Suppose we overload a Java method in a subclass and change the method in the following ways. Which may cause problems at run time if permitted by the type checker? Which would be okay? Explain in a sentence or two, appealing to basic subtyping principles.
- We change the method's parameter from `String` to `Object` in the subclass.

- We change the method's return type from `String` to `Object` in the subclass.
- We change the method's visibility from `private` to `public` in the subclass.
- We change the method's visibility from `public` to `private` in the subclass.

6. Suppose we add interfaces to IC. An interface is declared as illustrated below:

```
interface Moveable {
    void move(int dx, int dy);
}

interface Resizable {
    void resize(int dx, int dy);
}
```

Classes can implement one or more interface:

```
class Shape { }

class Rectangle extends Shape implements Moveable, Resizable {
    void move(int dx, int dy) { ... }
    void resize(int dx, int dy) { ... }
}

class Circle extends Shape implements Moveable, Resizable {
    void move(int dx, int dy) { ... }
    void resize(int dx, int dy) { ... }
}
```

We can then declare variables to have interface types in the usual way:

```
Rectangle r = new Rectangle();
Moveable m = r;
m.move(10,10);
...
Resizable rs = r;
rs.resize(-10,2);
```

One interface can also extend another interface, in which case the subinterface “inherits” all methods listed in the superinterface:

```
interface HideableMoveable extends Moveable {
    void hide();
}
...
HideableMoveable hm = ...;
hm.move(10,10);
hm.hide();
```

Describe the semantic checks you would need to add to your IC compiler to ensure that interfaces are used correctly. In particular:

- Extend the subtyping rules on page 3 of the IC specification to include interfaces. You may use the letters I, J, and K to denote interfaces, so that you can distinguish interface names from class names.

- (b) Describe the semantic checks one must perform for each interface declaration and each class that implements an interface. How would you need to change the `ast` and `syntab` packages to support interfaces? (A few sentences is sufficient.)
- (c) Consider the ternary operator again. Suppose we have the following declarations:

```
boolean b;
Shape s;
Rectangle r;
Circle c;
Moveable m;
Resizable z;
```

If we used the following in a program without typechecking them, which could lead to a type error at run time?

- i. `s = b ? s : r`
- ii. `m = b ? r : m`
- iii. `c = b ? r : c`
- iv. `s = b ? r : c`
- v. `m = b ? s : c`

- (d) Do your rules for ternary expressions from the previous question enable you to check each of these assignments properly? Are there other assignment expressions that your rules could not handle as well? If so, what are the issues and how can you handle them? (A few sentences is sufficient.)

7. Consider a C-like language that manipulates pointers. Statements and expressions have the following syntax:

$$\begin{aligned}
 e &\rightarrow n \mid x \mid \&x \mid *e \\
 s &\rightarrow x = e \mid x = \text{malloc}() \mid *x = e
 \end{aligned}$$

where n is an integer constant, x is a variable, and `malloc()` allocates an integer or a pointer on the heap (according to the declared type of x), and then returns a pointer to that piece of data. The only types are pointers and integers, but pointers can be multi-level pointers. The syntax for types is:

$$T \rightarrow \text{int} \mid T*$$

- (a) Write typing rules for all of the expressions and assignment statements. Use judgments of the form $E \vdash s$ for statements, and judgments of the form $E \vdash e : T$ for expressions.
- (b) Now let's extend the types in this language with two type qualifiers `taint` and `trust`. Tainted data represents data that the program received from external, untrusted sources, such reading from the standard input or reading from a network socket. All of the other data is `trusted`. Perl and other languages use tainting to, for example, prevent certain forms of security attacks on web scripts.

To model tainting, we extend the set of statements with a `read()` statement that reads an untrusted integer value from an external source:

$$e \rightarrow \dots \mid \text{read}()$$

The syntax for qualified types is:

$$\begin{aligned}
 T &\rightarrow Q R \\
 R &\rightarrow \text{int} \mid T * \\
 Q &\rightarrow \text{taint} \mid \text{trust}
 \end{aligned}$$

For instance, `trust ((taint int) *)` represents a trusted pointer to a tainted location, and `taint ((taint int) *)` denotes a tainted pointer to a tainted location.

Write appropriate typing rules for expressions n , x , $\&x$, $*e$, and `read()` for programs with qualified types. Also write a rule for `malloc`.

- (c) We want to prohibit the flow of values from untrusted sources into trusted portions of the memory. However, we want to allow flows of values from trusted locations to tainted locations. We can achieve this by defining an appropriate subtyping relation \leq between qualified types. First, we define an ordering between qualifiers:

$$Q \leq Q' \text{ iff } Q = \mathbf{trust} \text{ or } Q = Q'$$

We then use the subtyping rule:

$$\frac{[\text{SUBTYPE}] \quad Q \leq Q'}{QR \leq Q'R}$$

along with the standard assignment rule in the presence of subtyping:

$$\frac{[\text{ASSIGN}] \quad E \vdash x : T \quad E \vdash e : T' \quad T' \leq T}{E \vdash x = e}$$

to enforce the desired control over trusted values. For instance, these rules would make it possible to type-check this code fragment:

```

taint int x;
trust ((trust int) *) y;
y = malloc();
x = *y;

```

Prove that the above program type-checks by showing the proof trees for each of the two assignments.

- (d) Write the remaining rule for indirect assignments $*x = e$. Illustrate the use of this rule on a small program.
- (e) Consider the following, more general subtyping rules:

$$\frac{[\text{SUBTYPE 1}] \quad Q \leq Q'}{Q\mathbf{int} \leq Q'\mathbf{int}} \qquad \frac{[\text{SUBTYPE 2}] \quad Q \leq Q' \quad T \leq T'}{Q(T*) \leq Q'(T')}$$

Are these rules sound? If yes, argue why. If not, show a program fragment that type-checks, but yields a type error at run time.

8. Read sections I – III of the Sabelfeld and Myers paper, which covers the general issue of security and information flow and discusses a number of current research issues regarding how to ensure that confidential information does not accidentally leak out of a computation.

- What does non-interference mean in a security setting and how is it defined?
- What are the key ideas behind the type system of Section III?
- Show how to type check the valid and programs at the end of III-B, and explain why the invalid ones fail to check.
- What do the authors identify as the open issues in this area of research?

9. Read sections 1 – 3.1 of the Sampson, *et al.*, paper. The accompanying general-audience article can help give more context for this work.

- What is the main problem the EnerJ type system aims to solve?
- How does the type system, with *approximate* and *precise* types, relate to information flow or our trusted/tainted type system?
- EnerJ's endorsements are somewhat like casting. Does their treatment in EnerJ seem closer to casting in C (completely unchecked) or in Java (runtime object type checked against cast type)?